

The CFC and the CSD

How has the CSD evolved and what can the CFC expect next?

Suzanna Ward – Head of Data and Community

CFC15: March 2024



The Cambridge Crystallographic Data Centre



"The database was established in 1965 to fulfil a dream of myself and a great scientist, the polymath J.D. Bernal. We had a passionate belief that the collective use of data would lead to the discovery of new knowledge which transcends the results of individual experiments."

Olga Kennard , 1997 "From Private Data to Public Knowledge"



Polymath J.D. Bernal is typically credited for inspiring the original vision of the CSD #3, first Original Contents



Dr Olga Kennard, CCDC Founder



J.D. Bernal and group, including Olga Stonehenge, 1948



A wealth of information

- >1 million structures
 - > 100M 3D coordinates
- > 28 million bond lengths
 - > 2M unique distributions
- > 40 million valence angles
 - > 3M unique distributions
- > 14 million torsion angles
 - > 800K unique distributions
- > 2 million rings
 - > 400K unique distributions
- > 2 million hydrogen bonds
 - >30 million Isostar contacts





Images and graphics created using the CSD Python API and Flourish



15th CfC Meeting Series

15th CfC Meeting Series

The Cambridge Structural Database





5

The CSD and the CFC in 15 years





New elements during the CFC years





15th CfC Meeting Series

Helium added 2013 YEMTUH



Neon added 2016 AQUCOG



2023

HIRCUK

15th CfC Meeting Series

Changing chemistry





Changing nature of substances

15th CfC Meeting Series

8





*Katerina Vriza, University of Liverpool, PhD on Data driven discovery of functional molecular co-crystal

Changing chemistry & crystallography





Changing people and places





10

15th CfC Meeting Series

Number of CSD Communications addedCumulative number of entries



Increasing data in the CSD

15th CfC Meeting Series

The CSD contains every published structure including:

- Inc. ASAP & early view
- University repositories
- Historical articles
- Patents
- Thesis
- CSD Communications

Paid summer placements.

- Established >5 years ago
- 4 paid summer students a year
- Focussed on helping to improve wealth and value of data in the CSD
- Last 15 years of the CFC
- >500 patents
- >350 thesis publications





12

Valuable data from CFC members



Chart created by searching for published structures in the CSD that have been deposited using a company email address or are associated to crystallographer details with a company email address



>1,240 structures

• from past and present CFC members since 2008



Adapted from poster compiled and produced by the Njardarson Group (The University of Arizona) *Nature Reviews Drug Discovery 20, 85-90 (2021)

The CSD Drug and Pesticide Subsets



- Approved drugs in **ORUGBANK**
 - >15K in CSD Drug subset
 - +10K during 15 years of the CFC
 - >2K in Single-component subset
- In Pesticide Property DataBase
 - >1K in Pesticide subset



- Identified as molecules of interest in fight against COVID-19
 - >340 in COVID subset



iew Databases Results Help	_			
Entries in CSD version 5.44 (April 2023)	Results			Entry examples
Subsets in CSD version 5.44 (April 2023) 🕨	Best representative lists 🔸	<u> </u>		Note that this sorget makes use of functionality from the cookbook utility module.
Available Databases	CSD Drug subsets CSD MOF subsets ADPs available subset	CSD COVID-19 subset CSD Drug subset Single-component CSD Drug subset		The energy is and a proof of the star property without a processing water of the star proof of the sta
	CSD Pesticide subset Electron diffraction subset High pressure subset Hydrate subset Minimal disorder subset	Search Samp Samch Samch Samch Available Doublance: Deou Updates regularity C(2) varians 5.0((Bitvember 2021) = 1 spainte T	FRees Advanced Opposed 10 conclustes determined If Reduce P = 0.000 C = conclustes C = conclustes C = conclustes C = conclustes C = Other P = Non-disordered C = conclustes	The second section of the sec
	Polymorphic subset Significant disorder subset Teaching subset	You can search complete database()) or a subset (e.g., his found is a provision search) Select Subset Clear Subset Single query being used. Search will find structure: where this query is true In	No erros Not polymeric No ieres Conly of Single crystal structures C Powder structures C mly of Organics	
		Sert Search Cancel	C Organomatallic Reset	

(C)



Your data collections



L.N.Kalash et al. CrystEngComm 23 (2021) 5430-5442







Proprietary data curation

Data siloes and poor data quality can make finding, sharing and generating insights from data hard

Curation, management & consultancy towards FAIR industrial data

Helping to store experimental and predicted results together



Ensuring structural data is accessible across the organisation



15th CfC Meeting Series



18

Supporting your data management

- We are available to you to create, enhance, update and maintain your databases of in-house structures or to help you in this process as a service
- Service work has included:
 - Building databases
 - Auditing, updating and improving existing databases
 - Providing bespoke database training
 - Writing custom API scripts
 - Enhancing existing databases
 - Creating databases of CSP/DFT structures
 - Building bespoke subsets of data in the CSD
 - Data quality checks and enhancements







15th CfC Meeting Series

Where are we now and where next?

	Candidate	Planned	📕 In progress	Released		
Data	Capturing crystallisation conditions	New data fields available	Targeted bulk CSD improvements 2024	New data format for WebCSD/AS		Released
Data	Data versioning	Open licenses for deposited CIFs	Calculated semi-conductor data	Targeted bulk CSD improvements 2023		New data format for WebCSD/AS
	Evolving our knowledge bases	New data subsets		Improved mmCIF handling (Product)	F	
	Improved metadata for PXRD, ED, etc	Visualising disorder for CSD entries		InChIs via the API (Product)	L	Targeted bulk CSD improvements 2023
	Alternative ways to navigate/view entries	New data format for desktop release		Visualising Disorder in CIFs (Product)		Improved mmCIF handling (Product)
	New models of data					
	Improved stereochemistry representations					InChis via the API (Product)
Data integrity			• · · · · · · · · · · ·			Visualising Disorder in CIFs (Product)
Data integrity	New integrity checks in deposition	Data integrity filters	Additional data integrity metadata	Extended integrity checks for CSD entries	F	
Data partners		Improved and extended links	Further publisher workflows	Extended publisher workflows	L	Extended integrity checks for CSD entries
		CCDC visualisation in more journals		CSD available through PDF-5+ (ICDD)		
Data standards			Community CSP data standards			Extended publisher workflows
			Schema for crystallisation conditions			CSD available through PDF-5+ (ICDD)
			CoreTrustSeal re-certified			





19

Data releases 2023-24

- November 2023 Full data update
 - 10,837 new structures (11,411 new entries).
 - Taking the **total size** of the CSD to **1,254,560 structures** (1,284,316 entries).
 - Included new entries and improvements to over 190K existing entries.

Coming soon!

- First quarterly data release for 2024
- Data update for on-site WebCSD



structure determined in 2023.



CSD Refcode: NIWYUR.

CSD Refcode: NITRUH.

20

15th CfC Meeting Series







21

Database and architecture evolution

Evolving the database at the heart of our software portfolio

- Distributed deployment via cloud based systems
- A unified platform and modern structure to facilitate searching
- Extended capacity to cater for increased data volumes
- Flexible and extendable to enable additional data types and features

New database format already deployed to our public WebCSD platform



www.ccdc.cam.ac.uk/solutions/csd-core/components/csd/database-evolution/



Improving existing entries

Targeted improvements allow improved integrity, consistency, discoverability and value of data

- During 2023 over 190K existing CSD entries have been **improved**, including improvements to:
 - Melting point, recrystallisation and bioactivity/source fields.
 - The addition of more oxidation states.
 - The labelling of radicals and polymorphs.
- New bulk editing capabilities have accelerated improvements



Keeping up with techniques

Enhancing Entries





Standardising text fields e.g. habit, colour and steps to improve the AI readiness of the CSD



Labelling polymorphs in CSD Drug Subset

15th CfC Meeting Series S LCU-104 (I41/a) - Mercury File Edit Selection Display Calculate CSD-Community CSD-Core CSD-Materials CSD-Theory CSD-Particle CSD-Discovery CSD Python API Help x > 0.6Open... Ctrl+O 0.6 > x > 0.4Show Labels for Carbon atoms Recent Files with Occupancy Suppression Disorder: A. All Sketch Molecule... Ctrl+K a b c a* b* c* x- x+ y- y+ z- z+ x-90 x+90 y-90 y+90 z-90 z+90 ← → ↓ ↑ zoom- zoom+ *x* < 0.4 SMILES to Molecule Ctrl+Alt+M Auto Edit Structure on Load Toggle the x = occupancySave As... Ctrl+S disorder options POV-Ray Image... Print in 3D... Exit Ctrl+Q 0.6740

Visualising disorder

It works in combination with other functionality, including Edit, viewing and selecting hydrogen bonds, display Voids.

CCDC

23

Data integrity



Online deposition process with embedded data integrity checks (checkCIF) and internal duplicate checks (automated and manual).



Majority of data (>90%) reflects data in associated peer reviewed scientific articles. Strong relationships with major publishers, referee request service.



Team of expert PhD level editors and scientists at CCDC, including a data integrity scientist, that support the curation of data into the CSD.



Workflows to compare deposited datasets with the CSD through comparisons and assessments to identify some instances of fraud



Strong connections with IUCr, individuals, some publishers, membership of COPE and use of external resources (Retraction Watch)



Clear retraction policy and workflow



15th CfC Meeting Series

November 2023 Update on Our Ongoing Investigations in to Structures Associated with a Pre-print on a Papermill in Crystallography

Our investigations by our Data Integrity team following the pre-print on a prolific papermill in crystallography began in April 2022. In May 2022 we provided an update to say that we had found 992 structures in the Cambridge Structural Database (CSD) linked to publications named there. At this point, we also added a note to all impacted structures in the CSD which read "This structure is currently under review following a 2022 study of a prolific papermill https://doi.org/10.21203/rs.3.rs-1537438/v1."

Although the pre-print claims that the publications were fabricated, it did not claim that the data was fraudulent. So, shortly after the pre-print was published, we started more extensive investigations and discussions with publishers. For publications to be retracted evidence is required; obtaining definitive proof can vary depending on the dataset and this can be a complex situation.

In April 2023 we provided an update to confirm 209 of the implicated structures had been retracted from the CSD following the retraction of 125 associated publications. These retractions left 783 implicated structures under investigation.

Working closely with publishers and our community, we are following COPE guidelines as our investigations progress. We continue to add editorial comments to entries to highlight information that may be relevant so users can select fit for purpose data. The CSD portfolio also enables researchers to critically evaluate the data in the context of the >1.25 million structures in the CSD. Our collaborations with publishers on investigations have also led to retractions in the scientific literature as well as the CSD. In addition, we work closely with the <u>LUCr</u> and other data repositories in this field.

November 2023 Status

CCDC

In our last data release of 2023, scheduled for later this year, a further 152 entries implicated by the pre-print will be retracted. These retractions are already visible to users accessing the CSD via our web platforms and takes the total number of retractions related to the papermill pre-print in the CSD to 361.

We have updated the comments for 39 structures related to the papermill pre-print to say that following extensive review by the publishers and the CCDC there are currently no concerns on the data. A further 36 structures related to the papermill pre-print have been updated where our investigations have identified the entry has almost identical reflection data with another entry in the CSD with a different structural formula.

Investigations are continuing for the remaining 556 structures implicated by the pre-print. Cases where reflection data has not yet been made available can be especially challenging to assess. We will continue to work with publishers, reviewers and depositors and are grateful for how the community has come together to tackle the issues so far.

Work to broaden the automated data integrity checks and processes, conversations with publishers on how we can support their peer review processes as well as collaborations with the community remain a key priority for the CCDC.

Keep up to date with further developments here.



15th CfC Meeting Series

Where are we now and where next?







15th CfC Meeting Series

Where are we now and where next?







15th CfC Meeting Series

Training courses



CCDC Virtual Workshops

16th April

How crystal structure affects particle behaviour

30th April

Advanced functionality for visualization and analysis of structures in Mercury

ires in Me

14th May In-depth comparison of polymorphic structures using Mercury



Register now or share with your co-workers



https://www.ccdc.cam.ac.uk/community/training-and-learning/csdu-modules/

CDC

The CCDC Virtual Workshops are a series

beginners and more experienced users of

of hands-on, guided training sessions,

where you learn how to use different

components of CSD software. These sessions are free and open to both

The format is 90 minutes and Show One, Try One, Explore More:

· Show One: A guided demo of the

software by the CCDC tutors.

• Try One: Hands-on examples for

you to try with CCDC tutors on

Explore More: Learning outcomes

recap, challenges, and quizzes.

Registration is now open -

just scan the QR code

WHAT?

the CSD Software.

hand to help.

HOW?

FREE VIRTUAL WORKSHOPS

APRIL-MAY 2024

16 April - 1 pm (BST)

How Crystal Structure Affects Particle Behaviour

30th April - 4 pm (BST)

Advanced Functionality for Visualization and

Analysis of Structures in

==____ 14th May - 10 am (BST)

In-Depth Comparison

of Polymorphic

Structures Using

Mercury

Mercury

SCAN ME





27



Tributes to Olga Kennard



Celebrating the Life and Legacy of Dr Olga Kennard

Sunday 17 March 2024, Spring 2024 ACS National Meeting, New Orleans Organizers: Ian Bruno, Judith Currano, Carmen Nitsche, Suzanna Ward,

Start	Speaker	Title	
8:00	Suzanna Ward	The life of Dr Olga Kennard: A visionary leader	Remote
8:15	Stephen Burley	Protein Data Bank: From Two Epidemics to the Global Pandemic to mRNA Vaccines and Paxlovid	Remote
8:35	Rivka Isaacson	Olga Kennard as an icon of interdisciplinarity	In-person
8:55	Dirk Trauner	The chemist and the architect	In-person
9:15	Lawrence Falvello	Olga Kennard as research mentor in the mid-1970's: A DNA fragment, dipeptides, an early use of disorder tools, and cannabidiol	In-person
9:35	John Rumble	Olga Kennard and the Cambridge Crystallographic Data Centre: An Inspiration for Modern Scientific Data Science: An Outsider's View	In-person
9:55	Break		
10:05	Jeremy Sanders	Olga Kennard: The CSD and wider contributions to the community	Remote
10:20	Carolyn Brock	Approximate symmetry in organic crystal structures having Z'>1	In-person
10:40	Ian Bruno	From vision through data to knowledge: The insights of Dr Olga Kennard	In-person
11:00	Rajni Bhardwaj	The integration of CCDC Tools into Pharmaceutical Solid form Development	Remote
11:20	Diane Dickie	Teaching Crystallography to Chemistry Students Using the Cambridge Structural Database	In-person
11:40	Jürgen Harter	The future of structural science: Building on Olga Kennard's enduring legacy	In-person



CRYSTAL GROWTH ADESIGN Geometric Analysis and DFT Study of 2,2'-Dipyridylamine-Stabilized irst-Row Transition-Metal Comp Gradual Changes in the Aromaticity in a Series of Hydroxypyridine. Carboxylic Acid Derivatives and Their Effect on Tautomerism and Crystal Packing CRYSTAL GROW TH & DESIGN Hydrogen vs Halogen-Bonded R₂²(8) Rings in Organic Crystal DESIGN The Advantages of Flexibility: The Role of Entropy in Crysta Structures Containing C–H…F Interactions ameron I. G. Wilson, Ian Plesniar, Heike Kuhn, Jeff A Q R XXX De Judter, Published in Manufacture of Stationard Action

15th CfC Meeting Series

28

Thank you and questions

I think that the great ocean of truth is still in front of us and that we will continue to discover new aspects of this truth.²⁰

Dr Olga Kennard, Founder of the CSD

Olga Kennard - Bernal's Vision: From Data to Insight - J.D.Bernal Lecture 1995







15th CfC Meeting Series